

Linear models for regression

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ the output will be a real value (regression)

$$D = \{(x_m, t_m)\}_{m=1}^N$$

t_m is now a real value
input-output pairs

- ① The simplest linear model for regression is a linear combination of the input variables

$$y(x; w) = w_0 + w_1 x_1 + \dots + w_d x_d = w^T x$$

parameters input

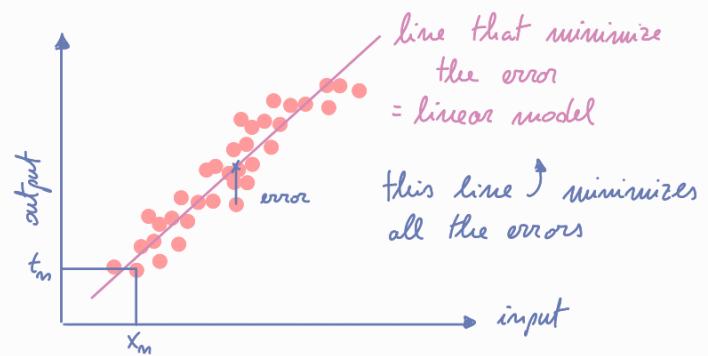
linear model of a linear function

$$x = [1, x_1, \dots, x_d]^T; w = [w_0, w_1, \dots, w_d]^T$$

This is known as linear regression.

The key property is that it is a linear function of the parameters. It is also, however, a linear function of the input and this imposes significant limitations on the model.

The goal is to find an approximation of the true function but the model is linear so it will find an hyperplane that fits the distribution of the dataset



Considering the previous limitations, we extend the class of models by using nonlinear functions of the input variables.

- ② Even more important than with classification is to consider basis functions and non linear transformations of the input.

We define a set of functions ϕ_i $i=0, \dots, n$ (in total they are $n+1$) and consider such model:

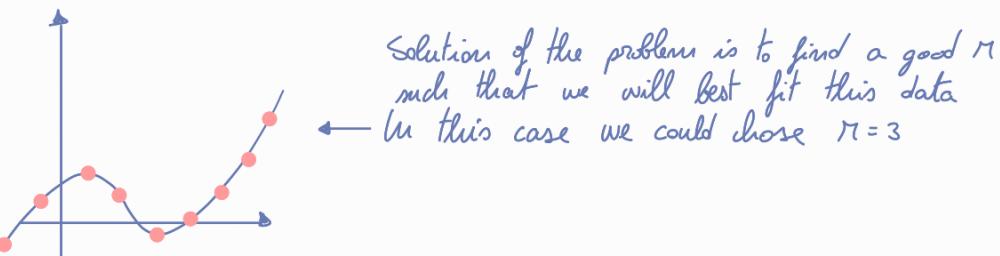
$$y(x; w) = w^T \phi(x) \quad \text{with } \phi(x) = \begin{bmatrix} \phi_0(x) \\ \vdots \\ \phi_n(x) \end{bmatrix} \quad \text{and } \phi_0(x) = 1 \quad \text{still } n \text{ degrees of freedom}$$

The resulting model is still linear in w but not in x (ϕ is usually non linear, otherwise this whole thing will be pointless).

e.g. we consider this particular transformation: $\phi = [x^0, x^1, \dots, x^M]^T$

$$y(x; w) = w^T \phi(x) = (w_0, w_1, \dots, w_M)(x^0, x^1, \dots, x^M)^T = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M$$

this is an M-degree polynomial



The choice of M is critical, it could lead to OVERFITTING: higher M could better fit the available data but we are interested in classifying new samples

There are several basis function → polynomial, radial, sigmoid, tanh

$$\phi = [x^0, x^1, \dots, x^n]^T$$

$$\phi_i(x) = e^{-\frac{(x-\mu_i)^2}{2s^2}}$$

$$\phi_i(x) = \sigma\left(\frac{x-\mu_i}{s}\right)$$

$$\sigma(a) = \frac{1}{1+e^{-a}}$$

$$\tanh(a) = 2\sigma(a) - 1$$

So far we just have seen an intuitive solution. What is the criterion to derive it?

Target values t in the dataset are given by the true function $y(x; w)$ affected by additive noise ϵ .

$$t = y(x; w) + \epsilon$$

We assume that the noise follows a gaussian distribution.

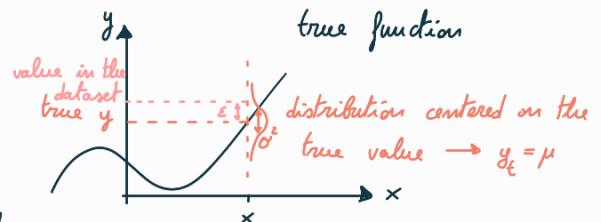
We just described the same problem from a probabilistic point of view.

$$P(t | x, w, \beta) = N(\epsilon | y(x; w), \beta^{-1})$$

likelihood

variance is described in terms of precision β ; precision is the inverse of the variance

the distribution is centered on the true value (zero-mean, we can translate it (?)



Now we can try to maximize the likelihood (minimize the negative log likelihood)

$$\text{maximize } P(\{t_1, \dots, t_N\} | x_1, \dots, x_N, w, \beta) = \prod_{n=1}^N N(t_n | w^\top \phi(x_n), \beta^{-1}) \quad \text{under the hypothesis that the samples are independent and identically distributed we can use the theorem of total probability}$$

or equivalently:

$$\begin{aligned} \text{minimize } \ln P(\{t_1, \dots, t_N\} | x_1, \dots, x_N, w, \beta) &= \sum_{n=1}^N \ln N(t_n | w^\top \phi(x_n), \beta^{-1}) \\ &= -\beta \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2 - \frac{N}{2} \ln (2\pi \beta^{-1}) \end{aligned}$$

$E_D(w)$
do not depend on w

$E_D(w)$ is what we need to minimize

Maximum likelihood (minimum log likelihood) under zero-mean Gaussian assumption correspond to least square error minimization

$$\underset{w}{\text{argmax}} \ P(\{t_1, \dots, t_N\} | x_1, \dots, x_N, w, \beta) \longleftrightarrow \underset{w}{\text{argmin}} \ E_D(w) = \underset{w}{\text{argmin}} \ \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2$$

sum of the squared errors:
sum over all the points in the dataset of the squared difference between the true value t_n and the prediction $w^\top \phi(x_n)$

linear regression can be solved with this optimization function that derives at the same time from the probabilistic formulation of the problem (maximum likelihood) and from the least square minimization (squared difference between the true value t_n and the prediction $w^\top \phi(x_n)$).

From the closed form solution:

- ① compute the gradient $E_D(w) = \frac{1}{2} (t - \phi w)^T (t - \phi w)$ square
 - with $t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix}$, $\phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_n(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_n(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_n(x_N) \end{bmatrix}$ $w = \begin{bmatrix} w_0 \\ \vdots \\ w_M \end{bmatrix}$ one weight for every element of the ϕ vector
 - ~~② $\nabla E_D(w) = 0 \Leftrightarrow \phi^T \phi w = \phi^T t$~~ using n instead of $n+1$ for simplicity
 - ③ isolate w from $\phi^T \phi w = \phi^T t$ $\phi^T \phi w = \phi^T t$
 $(M \times N)(N \times n) (n \times 1) = (M \times N)(N \times 1)$
 $(n \times n)(n \times 1)$ (M \times 1) $\alpha = 1$
- $w = (\phi^T \phi)^{-1} \phi^T t$ pseudo-inverse of $\phi \hat{=} \phi^+$

example

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$D = \{(x_m, t_m)\}_{m=1}^N$$

define t, w and ϕ

we assume this particular transformation:

$$\phi = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ & \vdots & & \\ 1 & x_N & x_N^2 & x_N^3 \end{bmatrix}$$

$\uparrow \quad \quad \quad \uparrow \quad \quad \quad \uparrow$
 $0 \quad \quad \quad \quad \quad \quad n=3$

the ϕ matrix will have N rows, one for each sample in the dataset and $n+1$ cols (number of features of the ϕ vector + ϕ_0)

$$t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix} \quad w = \begin{bmatrix} w_0 \\ \vdots \\ w_M \end{bmatrix}$$

our model will be:

$$y(x; w) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

$$(W_{ML}) = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{pmatrix} = \phi^+ t \quad \text{pinv}(\phi) \cdot t$$

w from maximum likelihood

We need an iterative algorithm

Why do we need an iterative algorithm if we have a close form solution?
 Because the ϕ matrix may be too large (N depends on the dataset size and M could also be very large)

Iterative algorithm : sequential learning

$$\hat{w} \leftarrow \hat{w} - p \nabla E_m$$

$$\hat{w} \leftarrow \hat{w} - p(t_m - w^T \phi(x_m)) \phi(x_m)$$

This algorithm converges for suitable small values of p ; its value needs to be chosen carefully.

Regularization

Regularization is a technique to control over-fitting. We add a regularization term to an error function so that the whole error function to be minimized takes the form:

$$\underset{w}{\operatorname{argmin}} \quad E_S(w) + \lambda E_W(w)$$

\uparrow
regularization coefficient

λ penalizes the weights that determine the overfitting that are too high in absolute value

A common choice for $E_W(w)$:

$$E_W(w) = \frac{1}{2} w^T w \rightarrow \text{this choice of regularizer is known as weight-decay because in sequential learning algorithms it encourages weight values to decay towards zero, unless supported by the data.}$$

Another choice:

$$E_W(w) = \sum_{j=0}^M |w_j|^q$$

Putting all together:

$$\underset{w}{\operatorname{argmin}} \quad \frac{1}{2} \sum_{m=1}^N (t_m - w^T \phi(x_m))^2 + \frac{\lambda}{2} w^T w$$

In general, when we add a parameter, we increase the flexibility of the model (better system that can more easily find a solution) but it needs to be well calibrated

Linear regression with regularization is called **Ridge Regression**

Everything can be extended to regression with multiple outputs

So far we have considered the case of a single target variable t . In some applications we may wish to predict $k > 1$ target variables which we collectively denote with the target vector t . This could be done by introducing a different set of basis functions for each component of t , leading to multiple, independent regression problems; a more interesting and more common approach is to use the same set of basis functions to model all the components of the target vector.

$$f: \mathbb{R}^d \rightarrow \mathbb{R}^k$$

$$\text{output } y = \begin{bmatrix} y_1 \\ \vdots \\ y_k \end{bmatrix} = \begin{bmatrix} -w_1^T \\ w_2^T \\ \vdots \\ w_k^T \end{bmatrix} \begin{bmatrix} | \\ \phi \end{bmatrix}$$

in each row we have the coefficients for a model

$$W^T (k \times n+1) \quad \left\{ \begin{array}{l} W \text{ is } (n+1 \times k) \\ W^T \text{ is } (k \times n+1) \end{array} \right.$$

This is still a linear model (linearity of the weight matrix).

The maximum likelihood is:

$$\ln P(T | X, W, \beta) = \sum_{m=1}^N \ln N(t_m | W^T \phi(x_m), \beta^{-1})$$

$$T = \begin{bmatrix} \vdots \\ -t_k \\ \vdots \\ 1 \dots k \end{bmatrix} \quad \begin{array}{l} \text{each row contains the output of the } k \text{ models;} \\ \text{it has one row for each sample in the dataset} \end{array}$$

Similarly as before we obtain:

$$W_{ML} = (\phi^T \phi)^{-1} \underbrace{\phi^T T}_{\text{pinv}(\phi)} \quad \phi^+$$